

Unsupervised Behavior Extraction via Random Intent Priors

Hao Hu · MIG · Tsinghua

10/9/2023



Machine Intelligence Group

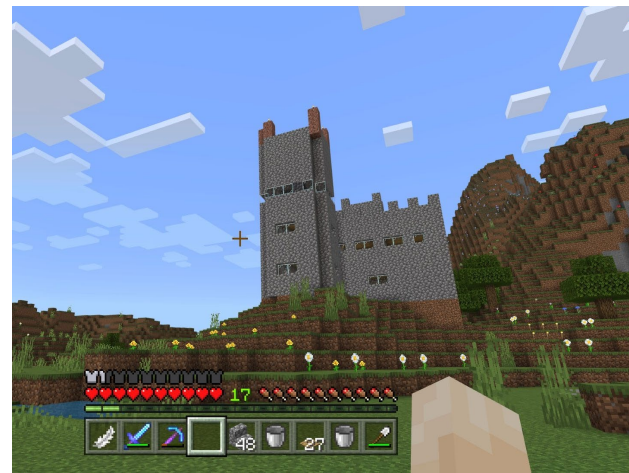


清华大学
Tsinghua University

交叉信息研究院
Institute for Interdisciplinary Information Sciences

Motivation

- Abundant reward-free data, containing useful human behaviors
- How to extract them effectively from offline data?



Random Neural Networks as Priors

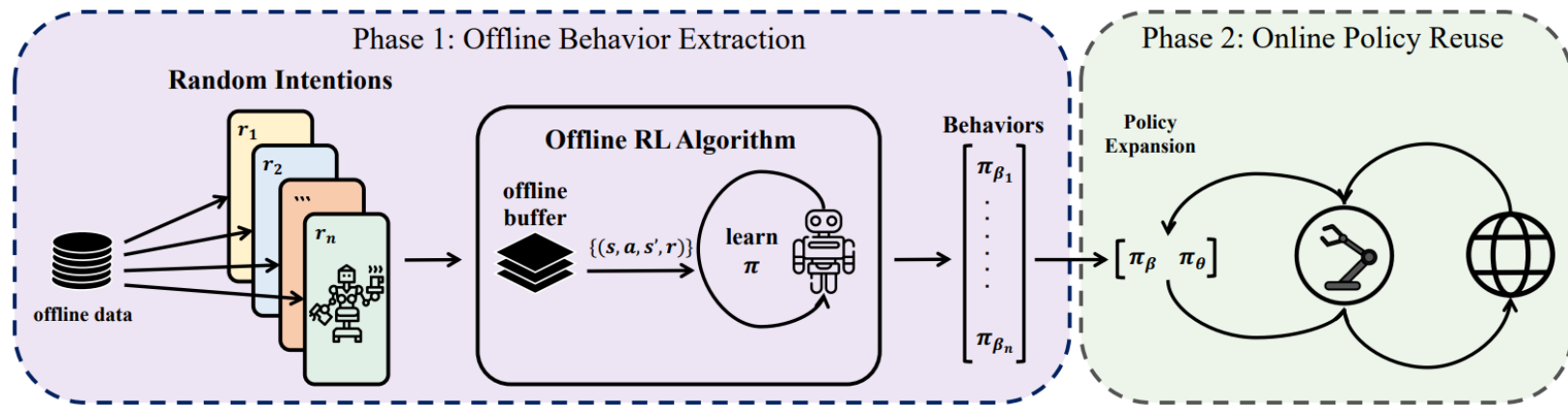


Figure 2: The framework of UBER. The procedure consists of two phases. In the first phase, we extract diverse and useful behaviors from the offline dataset with random rewards. In the second phase, we reuse previous behavior to accelerate online learning.



Policy Composition

- Policy set

$$\Pi = [\pi_\beta, \pi_\theta]$$

- Utility

$$P_{\mathbf{w}}[i] = \frac{\exp(Q_\phi(s, a_i)/\alpha)}{\sum_j \exp(Q_\phi(s, a_j)/\alpha)}, \quad \forall i \in [1, \dots, K]$$

- Composition

$$\tilde{\pi}(a|s) = [\delta_{a \sim \pi_\beta(s)}, \delta_{a \sim \pi_\theta(s)}] \mathbf{w}, \quad \mathbf{w} \sim P_{\mathbf{w}}$$



Algo: Random Neural Networks as Priors

Algorithm 1 Phase 1: Offline Behavior Extraction

- 1: **Require:** Behavior size N , offline reward-free dataset \mathcal{D}^{off} , prior intention distribution β
 - 2: Initialize parameters of N independent offline agents $\{Q_{\theta_i}, \pi_{\phi_i}\}_{i=1}^N$
 - 3: **for** $i = 1, \dots, N$ **do**
 - 4: Sample a reward function from priors $z_i \sim \beta$
 - 5: Reannotate \mathcal{D}^{off} as $\mathcal{D}_i^{\text{off}}$ with reward r_{z_i}
 - 6: **for** each training iteration **do**
 - 7: Sample a random minibatch $\{\tau_j\}_{j=1}^B \sim \mathcal{D}_i^{\text{off}}$
 - 8: Calculate $\mathcal{L}_{\text{critic}}^{\text{offline}}(\theta_i)$ and update θ_i
 - 9: Calculate $\mathcal{L}_{\text{actor}}^{\text{offline}}(\phi_i)$ and update ϕ_i
 - 10: **end for**
 - 11: **end for**
 - 12: **Return** $\{\pi_{\phi_i}\}_{i=1}^N$
-



Algo: Random Neural Networks as Priors

Algorithm 2 Phase 2: Online Policy Reuse

- 1: **Require:** $\{\pi_{\phi_i}\}_{i=1}^N$, offline dataset \mathcal{D}^{off}
 - 2: Initialize online agents Q_θ, π_w and replay buffer \mathcal{D}^{on}
 - 3: Construct expanded policy set $\tilde{\pi} = [\pi_{\phi_1}, \dots, \pi_{\phi_N}, \pi_w]$
 - 4: **for** each iteration **do**
 - 5: Obtain initial state from environment s_1
 - 6: **for** step $t = 1, \dots, T$ **do**
 - 7: Construct $P_{\tilde{\pi}}$ According to Equation (6)
 - 8: Pick an policy to act $\pi_t \sim P_{\tilde{\pi}}, a_t \sim \pi_t(\cdot | s_t)$
 - 9: Store transition (s_t, a_t, r_t, s_{t+1}) in \mathcal{D}^{on}
 - 10: **for** $i = 1, \dots, N$ **do**
 - 11: Calculate $\mathcal{L}_{\text{critic}}^{\text{online}}(\theta)$ and update θ
 - 12: **end for**
 - 13: Calculate $\mathcal{L}_{\text{actor}}^{\text{online}}(w)$ and update w
 - 14: **end for**
 - 15: **end for**
-



Theoretical Analysis: Completeness

Proposition A.1 (Proposition 4.1 restated). *For any behavior π , there exists an intention z with reward function r_z such that π is the optimal policy under r_z , i.e.*

$$V_{h,r_z}^\pi(\cdot) = V_{h,r_z}^*(\cdot). \quad (9)$$

Moreover, if π is deterministic, then π is the unique optimal policy under r_z , in the sense that $\forall \pi_{r_z}^, \pi_{r_z}^*(\cdot|s) = \pi(\cdot|s), \forall d^\pi(s) > 0$.*

Proof. Let $r_z(s, a) = \mathbb{1}(d_{s_0}^\pi(s, a) > 0)$, then we have

$$V_{h,r_z}^\pi(\cdot) = V_{h,\max} = V_{h,r_z}^*(\cdot). \quad (11)$$



Theoretical Analysis: Completeness

Theorem A.2 (Theorem 4.2 restated.). *Consider linear MDP as defined in Definition 2.1. With an offline dataset \mathcal{D} with size N , and the PEVI algorithm (Jin et al., 2021), the suboptimality of learning from an intention $z \in \mathcal{Z}$ with size $|\mathcal{Z}|$ satisfies*

$$\text{SubOpt}(\hat{\pi}; r_z) \leq 4c \sqrt{\frac{C_z^\dagger d^3 H^3 \iota}{N}}, \quad (12)$$

with probability $1 - \delta$ for sufficiently large N , where $\iota = \log \frac{4dHK|\mathcal{Z}|}{\delta}$ is a logarithmic factor, c is an absolute constant and

$$C_z^\dagger = \max_{h \in [H]} \sup_{\|x\|=1} \frac{x^\top \Sigma_{\pi_z, h} x}{x^\top \Sigma_{\rho_h} x},$$

with

$$\Sigma_{\pi_z, h} = \mathbb{E}_{(s,a) \sim d_{\pi_z, h}(s,a)} [\phi(s, a) \phi(s, a)^\top], \quad \Sigma_{\rho_h} = \mathbb{E}_{\rho_h} [\phi(s, a) \phi(s, a)^\top].$$



Theoretical Analysis: Coverage

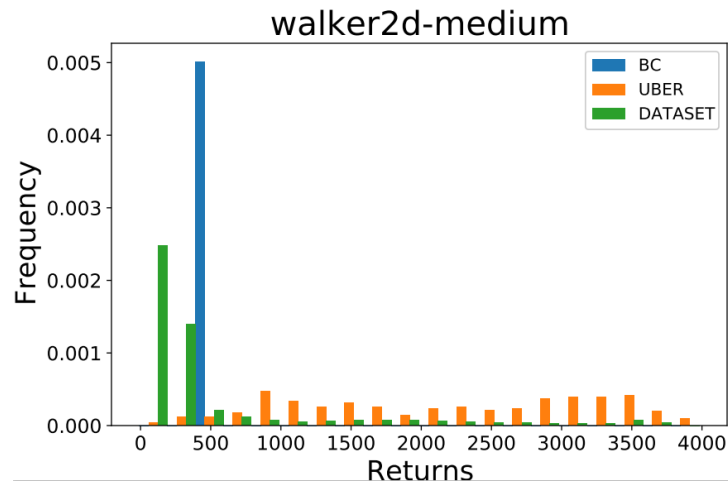
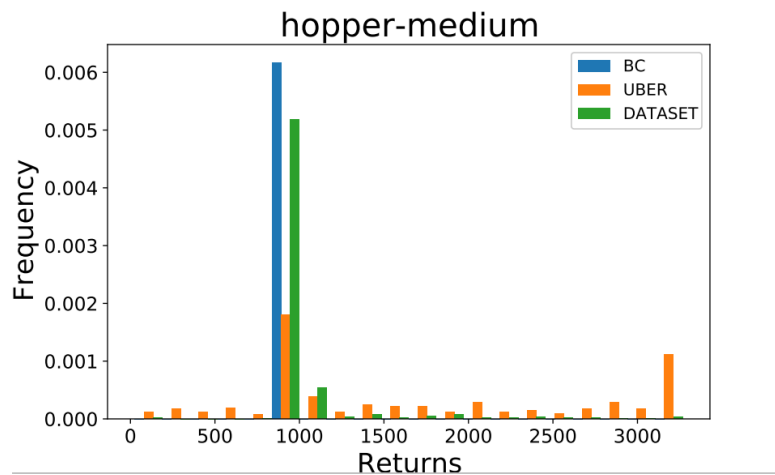
Theorem 4.3. *Assume the reward function admits a RKHS representation $\psi(s, a)$ with $|\psi(s, a)| \leq \kappa$ almost surely. Then with $N = c_0 \sqrt{M} \log(18\sqrt{M}\kappa^2/\delta)$ random reward functions, the linear combination of the set of random reward functions can approximate the true reward function with error*

$$\epsilon \leq c_1 \log^2(18/\delta)/\sqrt{M},$$

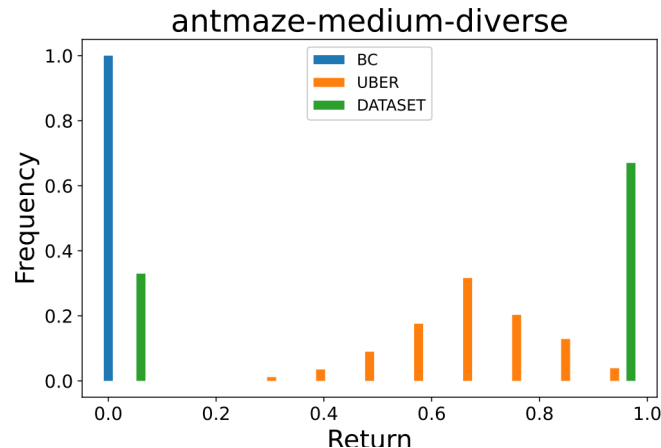
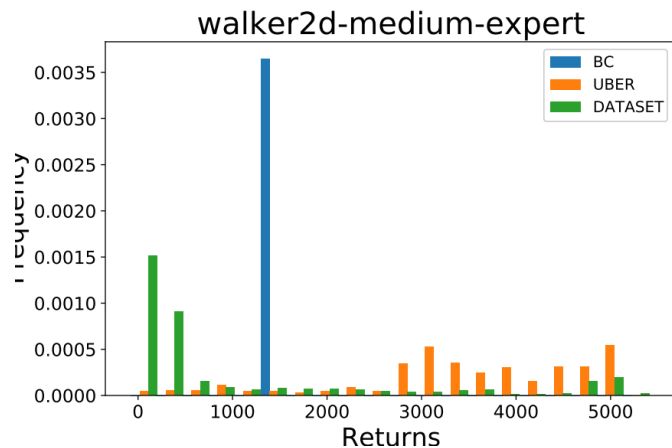
with probability $1 - \delta$.



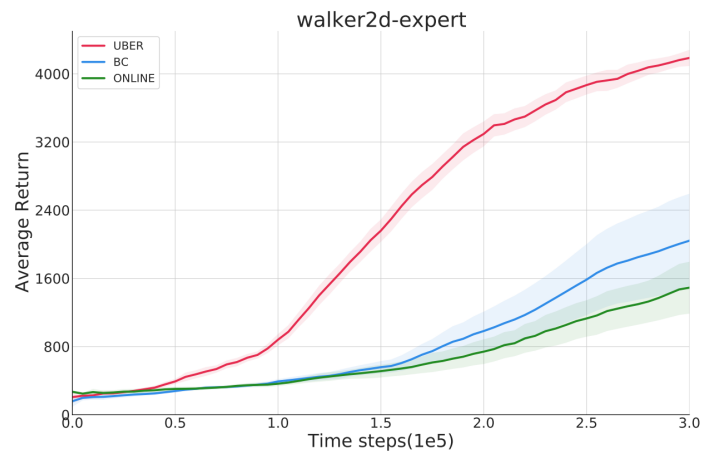
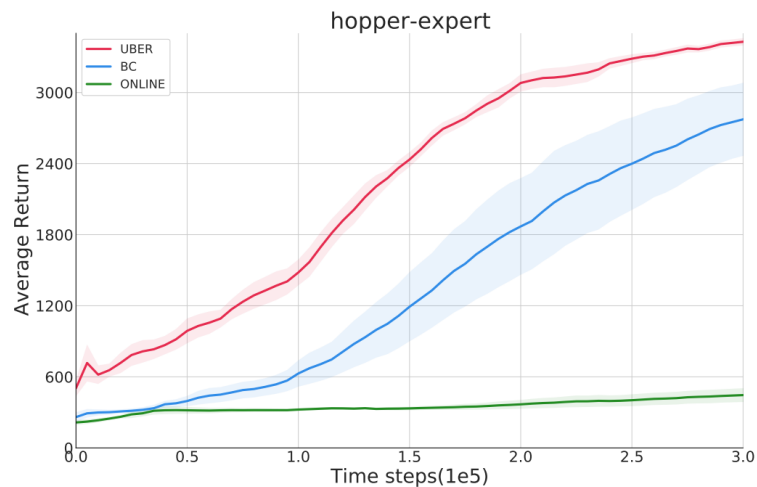
Experiments



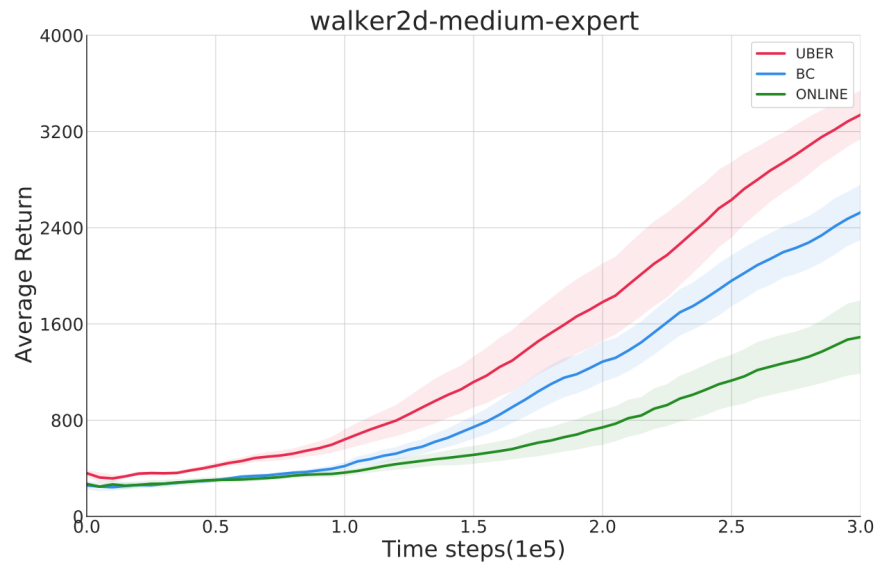
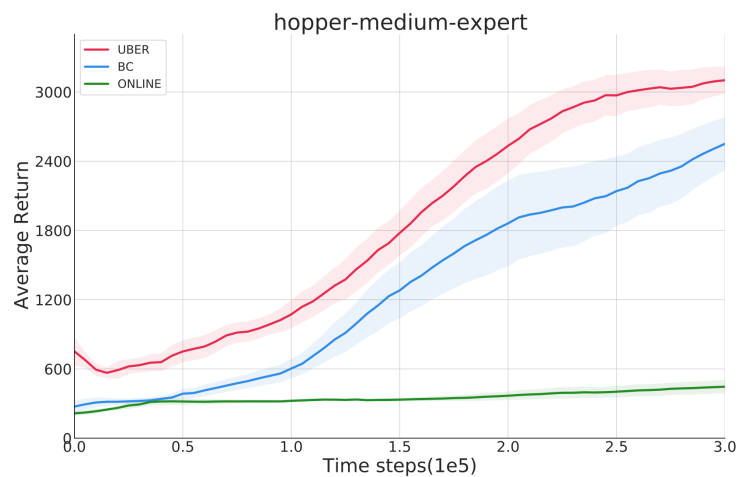
Experiments



Experiments

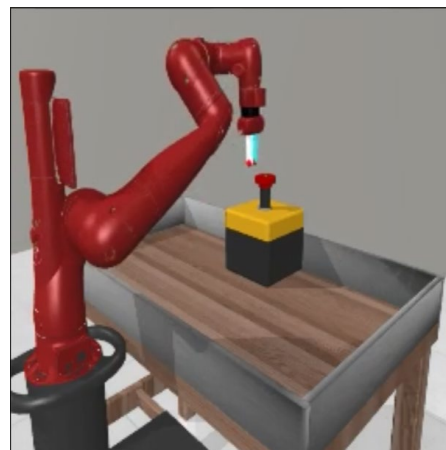
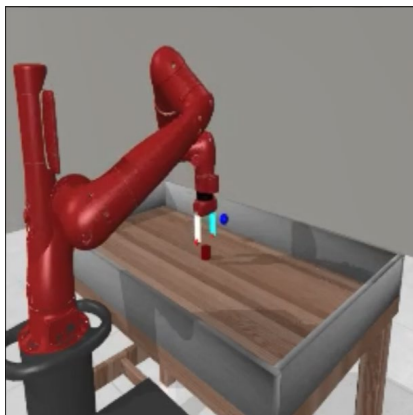


Experiments

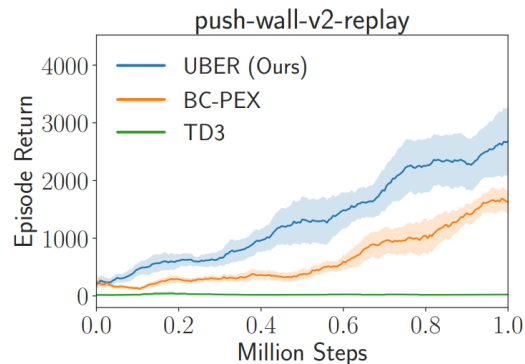
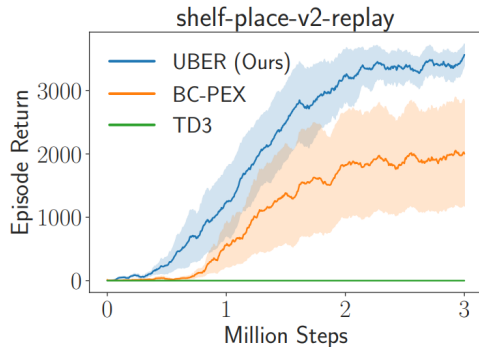
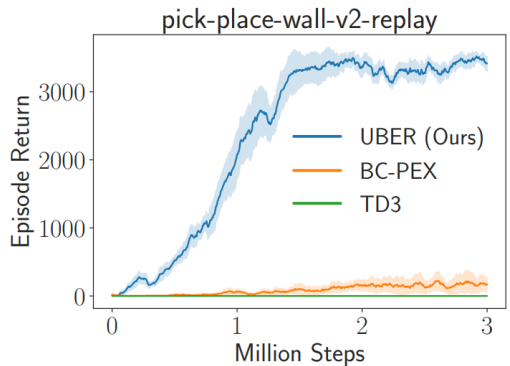
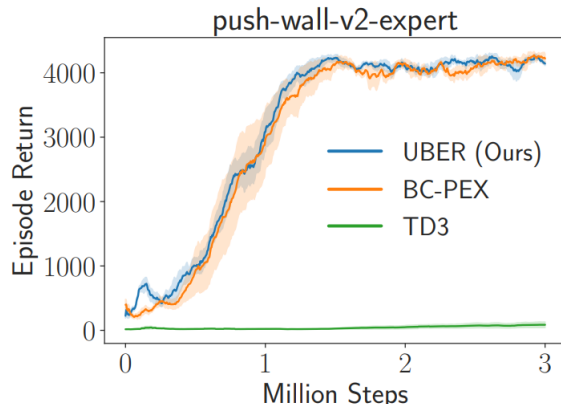
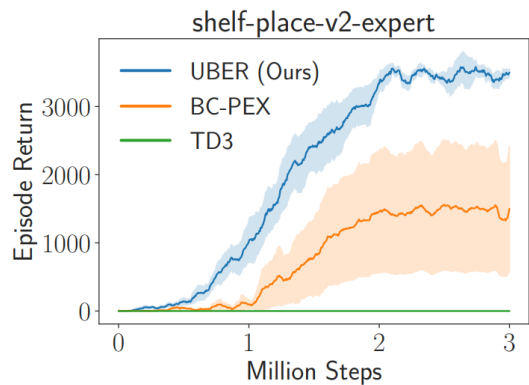
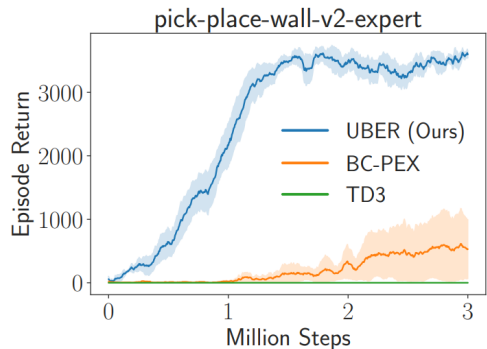


Multi-task: Meta-world

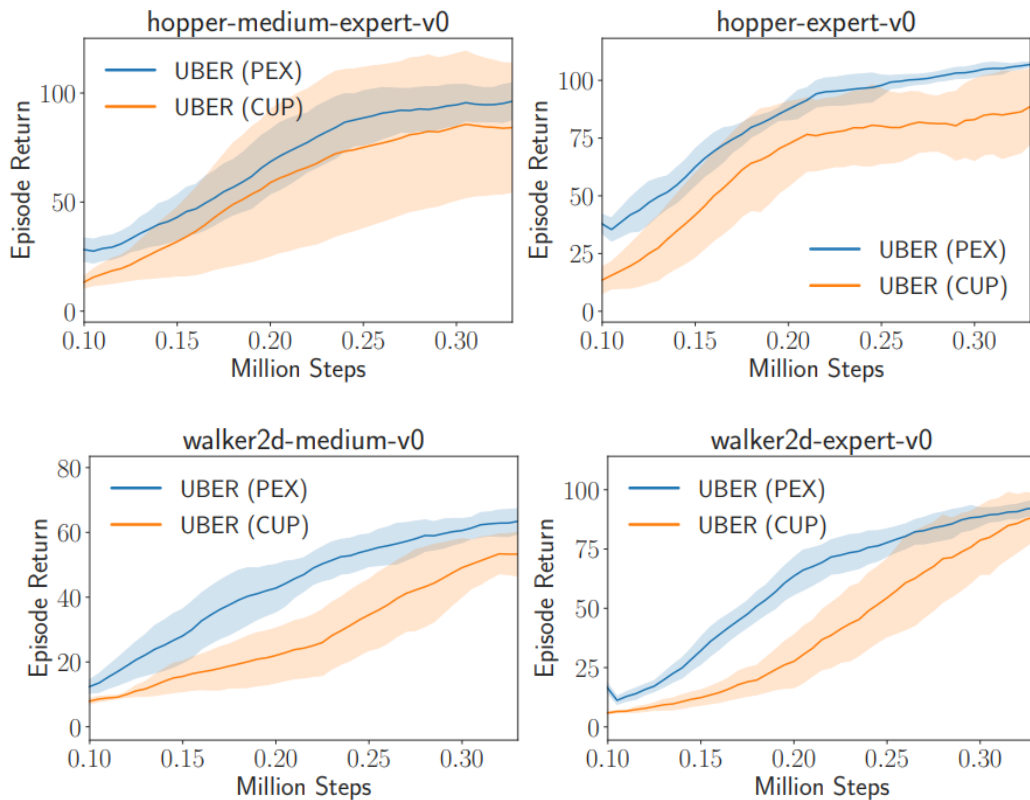
- Source: Push, Reach, Pick-place
- Target: Hammer, Peg-Insert-Side, Push-Wall, Pick-Place-Wall, Push-Back, Shelf-Place



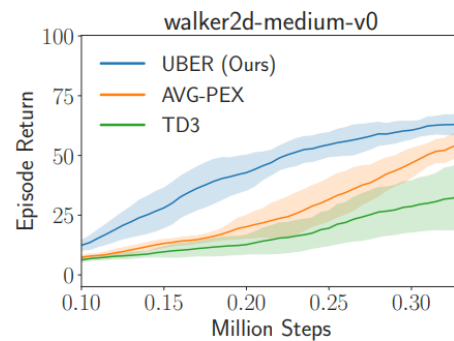
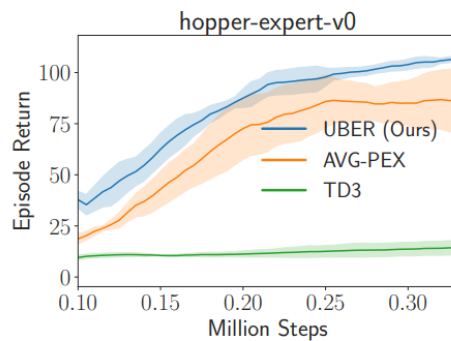
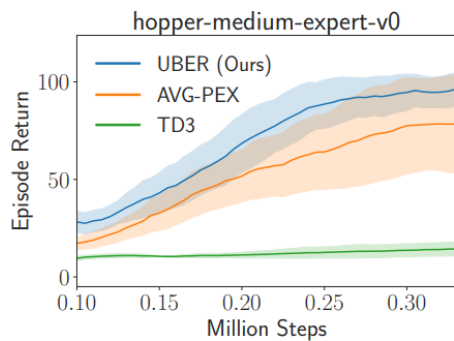
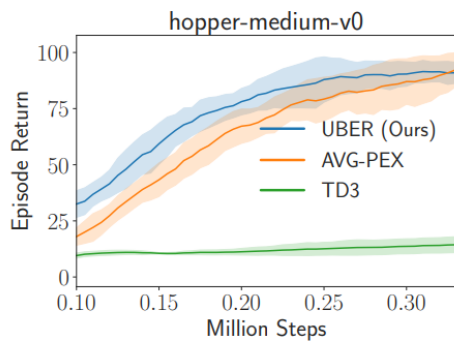
Results



Ablations



Ablations



Thanks!



Machine Intelligence Group



清华大学
Tsinghua University

交叉信息研究院
Institute for Interdisciplinary Information Sciences